

# Enhanced cloud/snow identification in snow mixed vegetation/soil areas based on machine learning techniques

Nan Chen<sup>1</sup>, Wei Li<sup>1</sup>, Charles Gatebe<sup>2</sup> and Knut Stamnes<sup>1</sup>

<sup>1</sup>Light and Life Laboratory, Stevens Institute of Technology, Hoboken, NJ, USA  
<sup>2</sup>NASA GSFC, Greenbelt MD

August 9, 2017



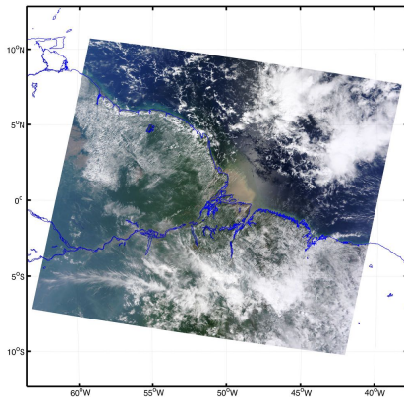
# Clouds and snow detection in remote sensing



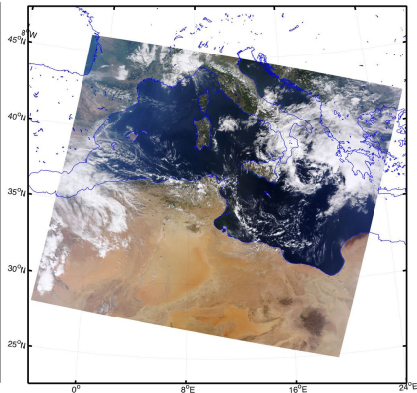
Earth image from NASA's Moderate-resolution imaging spectroradiometer (MODIS) sensor on July 11, 2005.

# Clouds over vegetation and other land types

South America

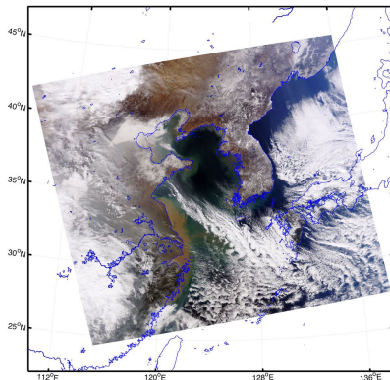


Sahara

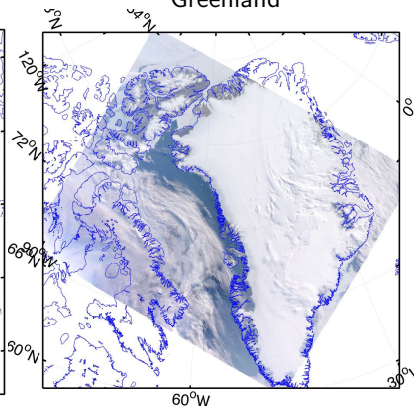


# Clouds over bright backgrounds (snow and ice)

East Asia



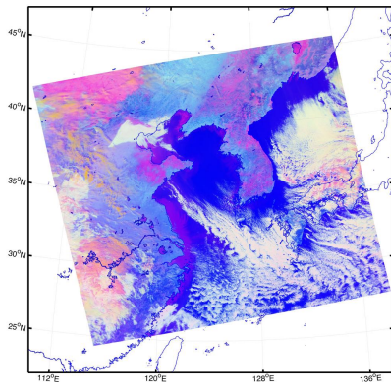
Greenland



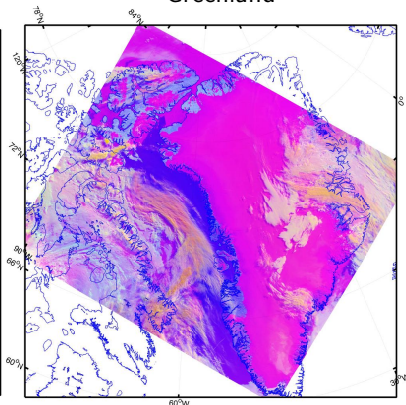


# Clouds over bright backgrounds (snow and ice)

East Asia



Greenland



False color RGB with R:  $0.645 \mu\text{m}$ ; G:  $2.13 \mu\text{m}$ ; B:  $10.8 \mu\text{m}$  brightness temperature.

# Methodology

- Combination of fixed or dynamic threshold tests.
- Use information from reflectance (VIS, NIR, SWIR) and thermal (TIR) channels.
- **Detect possible snow covered area** is important.

# NDSI based snow detection

Indexes for mapping snow cover using VIS and SWIR data were developed in the mid-1970s. The NDSI (Normalized Difference Snow Index) term was first coined by [Hall et al., 1995] and used to map snow using MODIS. Prior to that, [Dozier, 1987, Dozier, 1989] used a VIS/SWIR index algorithm to map snow using Landsat data. The NDSI is defined as:

$$NDSI = \frac{R_{VIS} - R_{SWIR}}{R_{VIS} + R_{SWIR}}$$

A fixed threshold (e.g.  $NDSI > 0.4$ ) is typically used to detect possible snow-covered pixels.

# Threshold tests - the MOD35 algorithm

Test No. #	Day Ocean	Night Ocean	Day Land	Night Land	Day Snow/Ice	Night Snow/Ice	Day Coast	Day Desert	Polar Day	Polar Night
RT <sub>11</sub> 13	✓	✓								
RT <sub>11,9</sub> 14	✓	✓	✓	✓	✓	✓	✓	✓		
RT <sub>1,1</sub> 15	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
R <sub>1,16</sub> 16	✓		✓		✓		✓	✓	✓	
RT <sub>1,1</sub> +RT <sub>1,2</sub> 17				✓		✓				✓
RT <sub>1,1</sub> +RT <sub>1,2</sub> 18	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
RT <sub>1,1</sub> +RT <sub>1,2</sub> 19	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
R <sub>1,16</sub> , R <sub>1,17</sub> 20	✓		✓				✓	✓		
R <sub>1,16</sub> 20								✓		
R <sub>1,16</sub> /R <sub>1,17</sub> 21	✓							✓		
RT <sub>1,1</sub> +RT <sub>1,2</sub> 22				✓		✓				✓
RT <sub>1,1</sub> +RT <sub>1,2</sub> 23	✓	✓								
So <sub>1</sub> , Temp <sub>1</sub> 24	✓	✓		✓						
RT <sub>1,1</sub> +RT <sub>1,2</sub> 25		✓								
RT <sub>1,1</sub> +Var 26		✓								

Threshold tests used in MOD35 algorithm  
 [Ackerman et al., 1998, Ackerman et al., 2010].

# Threshold tests - the ACCA algorithm

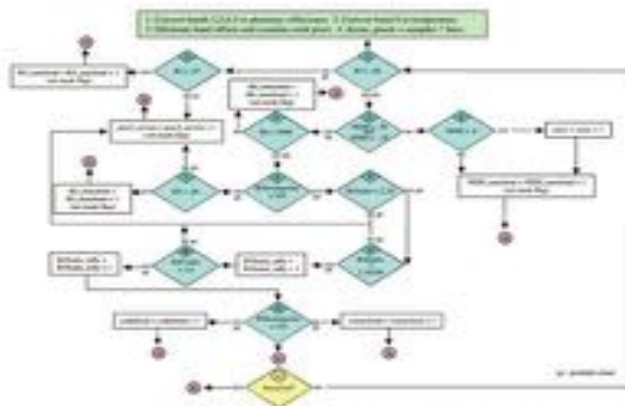


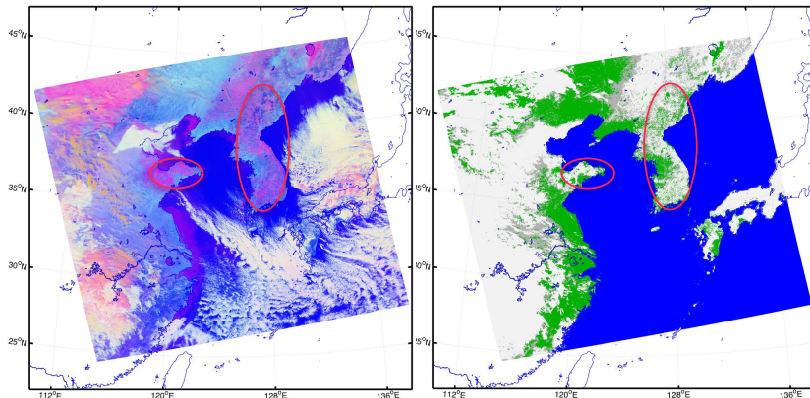
Figure 2. ACCA flow chart (1 of 4): spectral cloud identification.

Flowchart (1 of 4) of the Automated Cloud Cover Assessment (ACCA) algorithm [Irish et al., 2006] for the Landsat ETM+ sensor.

# Challenge: snow and snow mixed with vegetation areas

False color RGB image

MYD35 cloud mask



Aqua MODIS image over East Asia, Jan. 24, 2003. Color scheme of cloud mask figure: white/grey - clouds, blue - ocean, green - land. Only clouds detected over land areas are shown.

# Limitations of threshold based methods

- A large number of tests makes the logic very complicated and different tests may produce conflicting results.
- Fixed threshold settings make it difficult to handle complex surface mixing situations, especially when the surface is partially covered by snow.
- The accuracy of snow detection affects that of cloud detection because it determines the test combinations applied on each pixel.
- A relatively large number of satellite channels are needed to detect clouds as well as snow/ice (e.g. MOD35 uses 10 reflectance bands and 9 IR bands).
- Tests using thermal IR channels are not reliable over snow-covered areas, especially over Greenland and Antarctica due to frequent temperature inversions.

# Previous attempts using machine learning techniques

From the 1990s machine learning techniques (MLTs) are being used to detect clouds. Due to the complexity of the problem, **supervised machine learning techniques** are usually employed. These methods include:

- support vector machine (SVM)
- artificial neural network (ANN)
- logistic regression (LR)

**High quality training data sets** are the key to achieve high performance using machine learning algorithms.



# How to get training data?

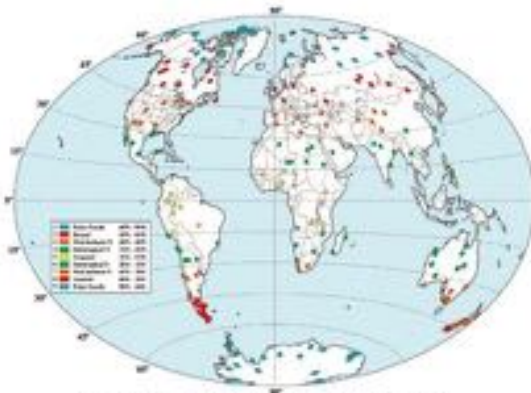


Plate 7. 500 random Landsat-4 locations in nine latitudinal zones.

The ACCA reference data set [Irish et al., 2006], one of the most commonly used validation/training data sets for cloud mask algorithms, is a **human identified data set**. It contains 212 scenes from 188 Landsat World Reference System (WRS) stations.

# Limitations of human identified training data set

- **Difficult to prepare** (need humans to identify millions of pixels).
- **Difficult to achieve full coverage** of solar/viewing geometries as well as different geolocations.
- **Difficult to account for seasonal land cover changes**, especially during the snow fall events.
- **Difficult to apply to a different sensor** due to sensor spectral characteristic changes.

## New Approach: Simulation + machine learning

Can we use a **simulated** TOA reflectance dataset to train machine learning algorithms?

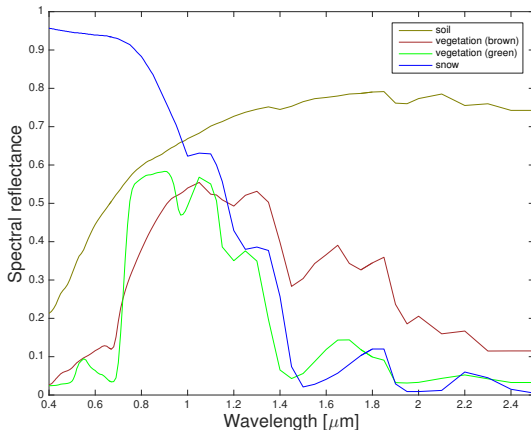
If possible, then we would have the following advantages over a human-identified dataset:

- There is no need for humans to identify hundreds of images with millions of pixels, which greatly saves human effort.
- The training dataset can cover the full range of possible solar/viewing geometries.
- Easy to apply to different sensors; only new training datasets are needed.

The biggest challenge in building a simulated dataset, is to **simulate complicated land surface types and account for various possible mixing conditions**.

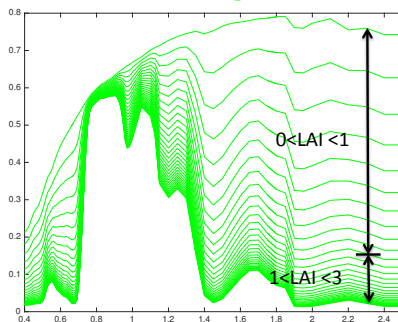
# How do we simulate complicated land reflectance?

Available land surface reflectance models, such as the Soil-Leaf-Canopy (SLC) model (Verhoef et al., 2007), can be used in **radiative transfer models (RTMs)** to simulate Top of Atmosphere (TOA) reflectances for an atmosphere overlying different land surfaces.

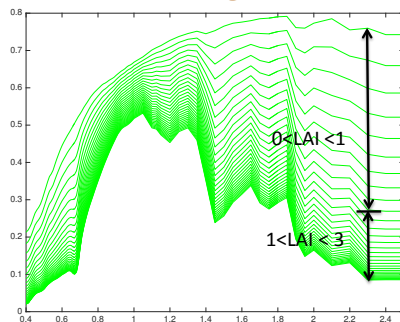


# Changing surface parameters

Green vegetation



Brown vegetation



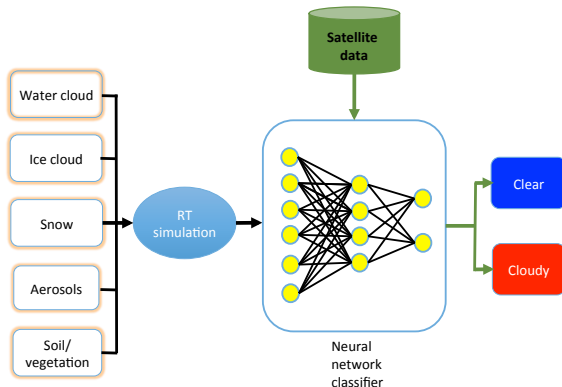
## Sub-pixel reflectance mixing

In order to better handle the case of fractional snow cover, we assume the following linear mixing rule for the reflectance of pixels with snow fraction,  $f$ :

$$R_{\text{mix}} = (1 - f) \times R_{\text{land}} + f \times R_{\text{snow}}.$$

By randomly changing the snow fraction  $f$  and different snow/land parameters, we can simulate the TOA reflectance of different snow-mixed-vegetation/soil cases.

# Neural network based cloud detection



- Over 10 million clear-sky and cloudy cases as the input.
- A simple classification neural network with one hidden layer.

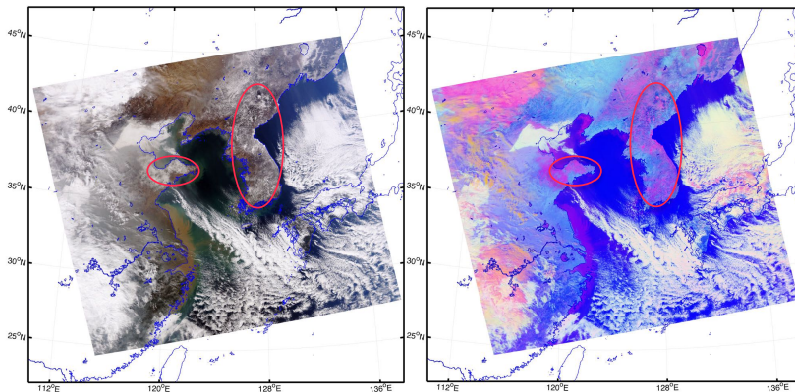
# Training for different sensors

The new neural network based algorithm can be configured (trained) for different numbers of channels. A 6-channel configuration was created and tested for Aqua MODIS and a 3-channel configuration for AVHRR-3.

<b>Sensor</b>	VIS chan ( $\mu\text{m}$ )	NIR chan ( $\mu\text{m}$ )	SWIR chan ( $\mu\text{m}$ )
3-channel	0.66	0.86	2.13
6-channel	0.47, 0.55, 0.66	0.86	1.24, 2.13

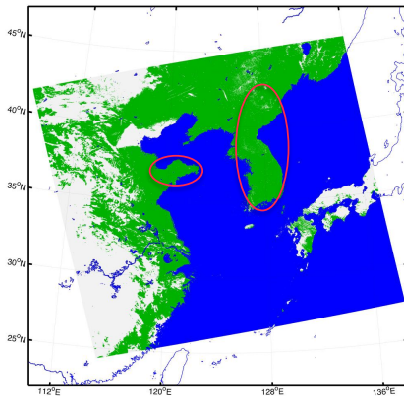


# Image based results: snow-mixed vegetation/soil area

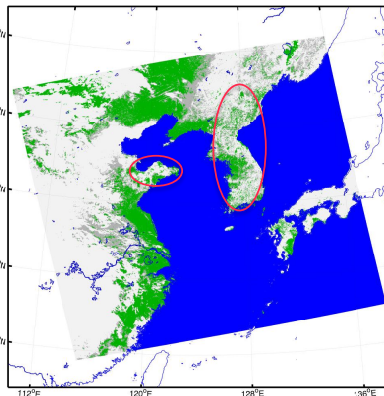


# Image based results: snow-mixed vegetation/soil area

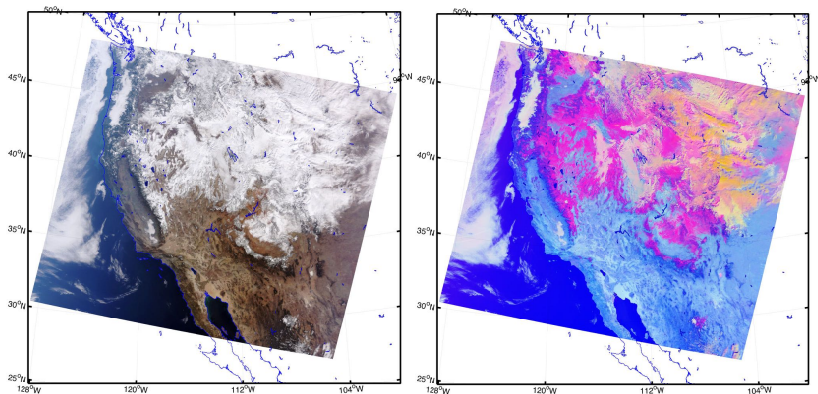
SCM cloud mask



MYD35 cloud mask

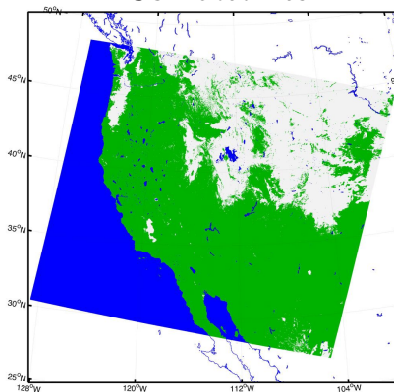


# Image based results: snow-mixed vegetation/soil area

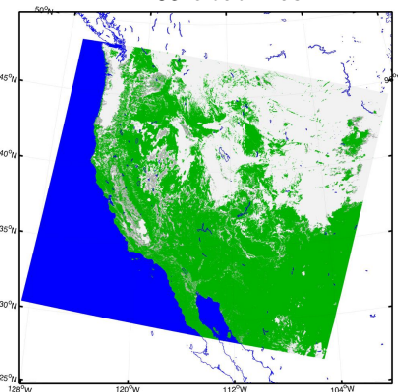


# Image based results: snow-mixed vegetation/soil area

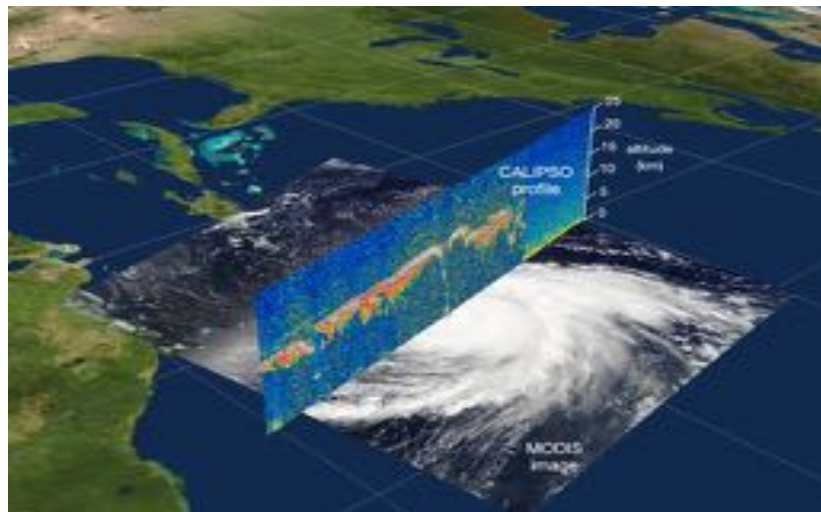
SCM cloud mask



MYD35 cloud mask



# CALIOP collocated with Aqua MODIS



# Validation using CALIOP observations as the benchmark

- Cloud-Aerosol Lidar with Orthogonal Polarization (**CALIOP**) is a lidar onboard the CALIPSO satellite that **provides high-resolution vertical profiles** of aerosols and clouds.
- Collocated Aqua MODIS/CALIOP data provide **the most reliable assessment** of cloud mask results by using CALIOP's active cloud detection scheme.
- MOD35 collection 6 product (MYD35 when using Aqua MODIS data) and CALIOP 1 km cloud layer product for the whole year of 2008 are used in the comparison.

# Test criteria in the comparison

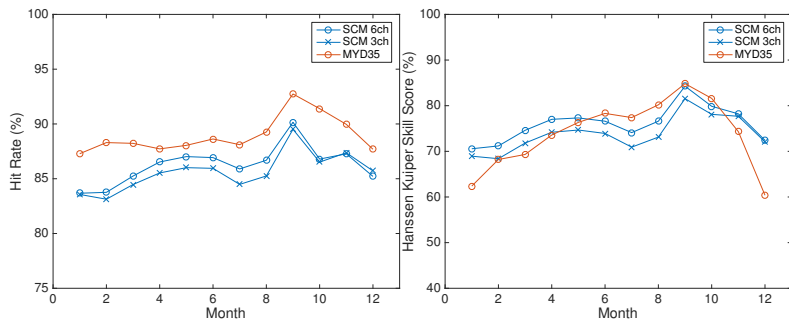
Hit Rate (HR)

$$HR = \frac{N_{cld,hit} + N_{clr,hit}}{N_{total}} \quad (1)$$

Hanssen-Kuipers Skill Score or **True Skill Score (TSS)**

$$TSS = \frac{(N_{cld,hit} \cdot N_{clr,hit} - N_{cld,miss} \cdot N_{clr,miss})}{(N_{cld,hit} + N_{cld,miss}) \cdot (N_{clr,hit} + N_{clr,miss})} \quad (2)$$

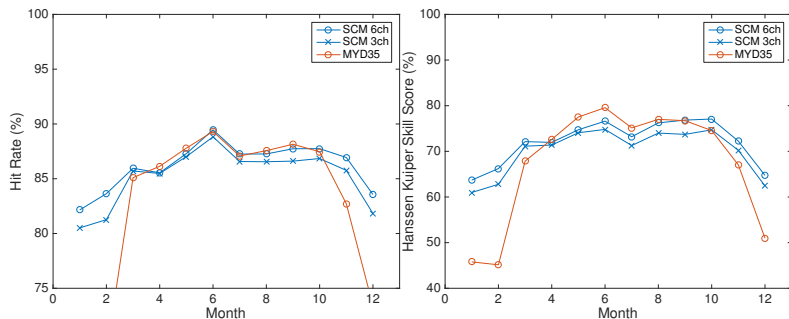
# Test against MOD35 over mid-latitude region (Europe)



HR and TSS of our algorithm (SCM) and MOD35 (MYD35 for Aqua MODIS) over Europe for the year 2008.



# Test against MOD35 over mid-latitude region (East Asia)



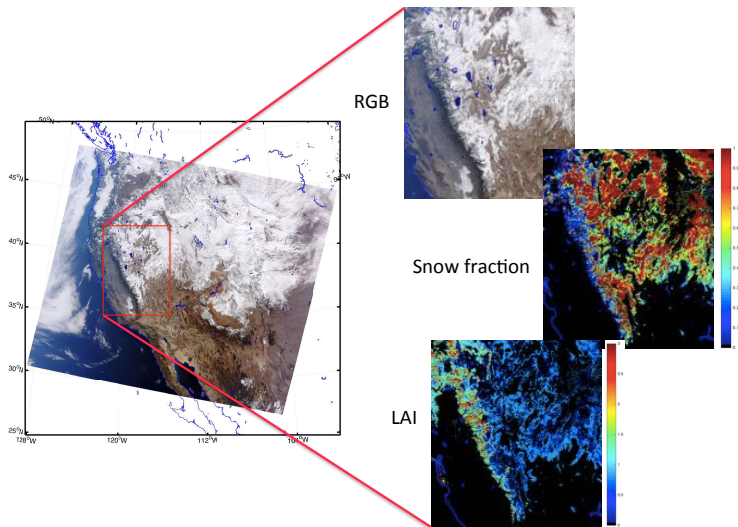
HR and TSS of our algorithm (SCM) and MOD35 (MYD35 for Aqua MODIS) over East Asia for the year 2008.

# Snow detection

How about snow detection? Can we **detect snow without thresholds**?

- Use the same dataset in cloud mask training (only the clear-sky cases).
- A separate neural network is trained for estimation of **snow fraction** or other parameters such as **Leaf Area Index (LAI)**.
- Different vegetation types (green/brown) and soil types are included in the training dataset.

# Results: snow fraction and LAI estimation



# Summary

- The new algorithm, without any thresholds, **consistently performs cloud detection using fewer channels** than threshold-based methods.
- Its performance is **significantly better than MOD35** (with higher TSS throughout) in the winter seasons when the surface is partially (or fully) covered by snow.
- It can **easily be re-configured for other sensors** (e.g. AVHRR and/or Landsat) with similar performance.
- It can **estimate the snow fraction** for each pixel, which goes beyond the traditional binary NDSI snow detection.
- Other parameters like **LAI or fAPAR can also be estimated** in a similar manner.



Thank you !!  
Questions ?

# Bibliography I



Ackerman, S., Frey, R., Strabala, K., Liu, Y., Gumley, L., and Baum, B. (2010).  
DISCRIMINATING CLEAR-SKY FROM CLOUD WITH MODIS ALGORITHM  
THEORETICAL BASIS DOCUMENT (MOD35).

Technical Report October, Cooperative Institute for Meteorological Satellite  
Studies, University of Wisconsin - Madison.



Ackerman, S. A., Strabala, K. I., Menzel, W. P., Frey, R. A., Moeller, C. C., and  
Gumley, L. E. (1998).

Discriminating clear sky from clouds with MODIS.

*Journal of Geophysical Research*, 103(D24):32141.



Dozier, J. (1987).

Remote sensing of snow characteristics in the southern Sierra Nevada.

*In Large Scale Effects of Seasonal Snow Cover*, number 166.



Dozier, J. (1989).

Spectral Signature of Alpine Snow Cover from the Landsat Thematic Mapper.

*Remote Sensing of Environment*, 22(February):9–22.

# Bibliography II



Hall, D. K., Riggs, G. A., and Salomonsont, V. V. (1995).

Development of Methods for Mapping Global Snow Cover Using Moderate Resolution Imaging Spectroradiometer Data.

*Remote Sensing of Environment*, 34:127–140.



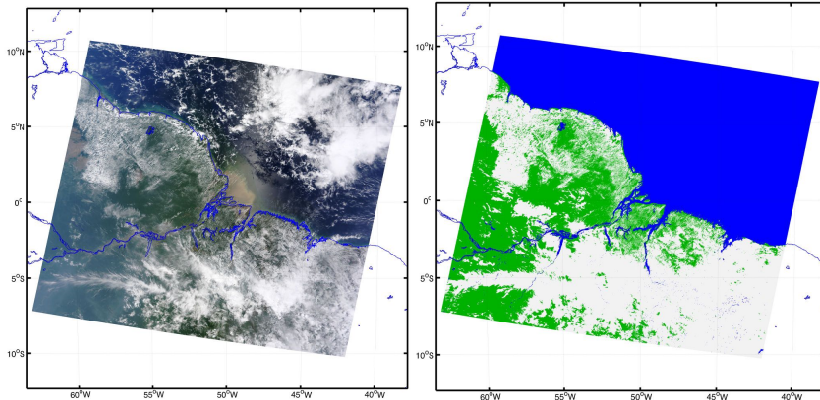
Irish, R. R., Barker, J. L., Goward, S. N., and Arvidson, T. (2006).

Characterization of the Landsat-7 ETM+ Automated Cloud-Cover Assessment (ACCA) Algorithm.

*Photogrammetric Engineering & Remote Sensing*, 72(10):1179–1188.

# Image based result: vegetated land area

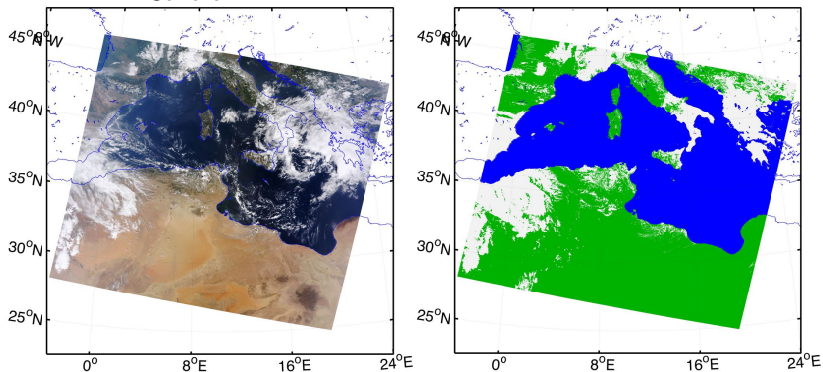
South America



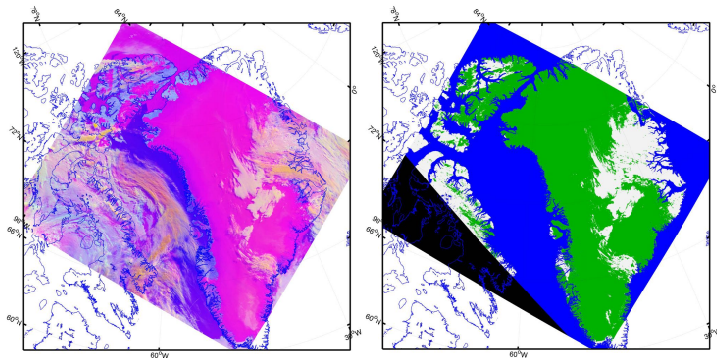


# Image based result: desert area

Sahara

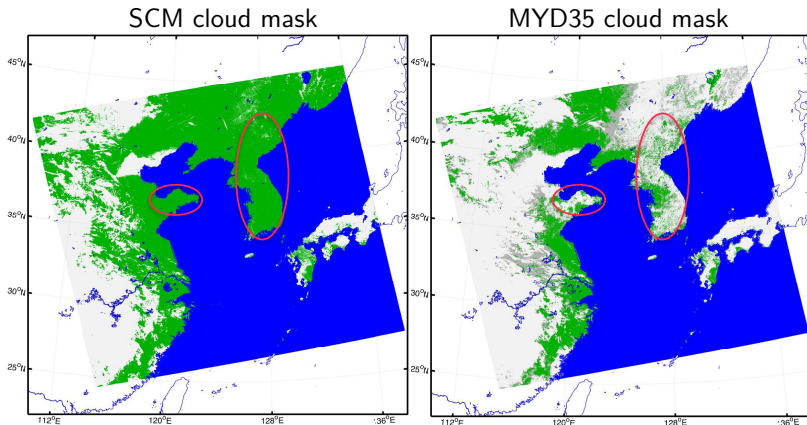


## Image based results: snow-covered area



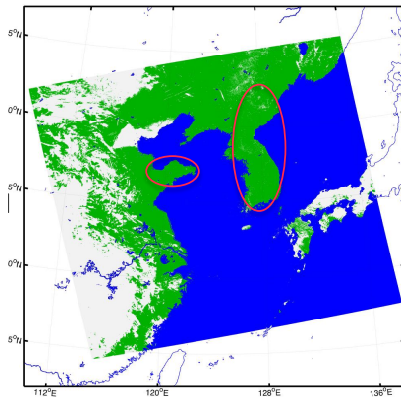
Comparison of cloud masks applied to an Aqua MODIS image over Greenland, Jul. 09, 2015. Left: False color RGB image; Right: SCM cloud mask. Color scheme of cloud mask figures: white/grey - clouds, blue - ocean, green - land.

## Image based results: snow-mixed vegetation/soil area



## Image based results: snow-mixed vegetation/soil area

SCM cloud mask



MYD09 cloud mask

